



DEPARTAMENTO
DE COMPUTACION

Facultad de Ciencias Exactas y Naturales - UBA

Trabajo Práctico II

Programación SIMD

Organización del Computador II
Segundo Cuatrimestre de 2014

Integrante	LU	Correo electrónico
Agustina Aldasoro	86/13	agusalaldasoro@gmail.com
Maximiliano Rey	37/13	rey.maximiliano@gmail.com
Ignacio Tirabasso	718/12	ignacio.tirabasso@gmail.com



Facultad de Ciencias Exactas y Naturales
Universidad de Buenos Aires

Ciudad Universitaria - (Pabellón I/Planta Baja)

Intendente Güiraldes 2160 - C1428EGA

Ciudad Autónoma de Buenos Aires - Rep. Argentina

Tel/Fax: (54 11) 4576-3359

<http://www.fcen.uba.ar>

Resumen

En el presente trabajo se describen los beneficios de la programación en Lenguaje Ensamblador bajo el modelo de programación SIMD mediante el uso de instrucciones SSE.

Índice

1. Objetivos generales	3
2. Contexto	3
3. Implementación en Assembler	4
3.1. Filtro CropFlip	4
3.2. Filtro Sierpinsky	6
3.3. Filtro Bandas	10
3.4. Filtro Motion Blur	14
4. Enunciado y solución	18
4.1. Mediciones	18
4.2. Filtro <i>cropflip</i>	18
4.3. Filtro <i>Sierpinski</i>	27
4.4. Filtro <i>Bandas</i>	29
4.5. Filtro <i>Motion Blur</i>	31
5. Conclusiones y trabajo futuro	32

1. Objetivos generales

El objetivo de este Trabajo Práctico es evaluar la eficiencia del modelo de programación SIMD mediante la implementación de diversos algoritmos en lenguaje Ensamblador utilizando instrucciones SSE.

Las mediciones se realizan mediante pruebas empíricas del código frente a algoritmos que cumplen la misma especificación, implementados en un lenguaje de alto nivel (C).

En este proyecto, los algoritmos a implementar se basaron en el procesamiento de imágenes y video, en el cual el uso del modelo SIMD es provechoso.

2. Contexto

Abordando el objetivo de este trabajo, realizamos una experimentación enfocada en reducir el tiempo de cómputo de un programa basado en dos sucesos que tienen la capacidad de afectarlo.

Por un lado se encuentra la *capacidad de cómputo*, la cual limita la cantidad de operaciones aritméticas que el procesador puede paralelizar. Si el programa ejecuta un uso intensivo en operaciones aritméticas, al añadirle nuevas operaciones de esta índole se van a necesitar más ciclos de clock para ejecutarlas.

El otro cuello de botella importante es el *ancho de banda de la memoria*. Cuando el programa ejecuta una instrucción que implica un acceso a memoria, se precisan más ciclos de clock para que la memoria responda, en particular si el dato no se encuentra en el cache.

Diseñamos nuestros casos de testeo con el fin de observar para cada programa cuál es el factor que determina el tiempo de cómputo.

Nuestra experimentación se centra en la cantidad de ciclos de clock que transcurren desde el inicio hasta el final de la ejecución del programa.

3. Implementación en Assembler

3.1. Filtro CropFlip

Luego de pushear los cinco registros a utilizar, almacenamos:

```
Pushear la base de la pila y los registros a utilizar.
Alinear la pila.
mov r12d, [rbp+16] ;tamx
mov r13d, [rbp+24] ;tamy
mov r14d, [rbp+32] ;offsetx
mov r15d, [rbp+40] ;offsety
Limpiamos la parte alta de estos registros haciendo un mov rXd, rXd.
```

Utilizamos el registro **r10** como el *y-actual* (fuente) y el registro **r11** como el *x-actual* (fuente). Recorremos la imagen fuente desde arriba hacia abajo, de izquierda a derecha.

```
mov r10,r15 ; y
mov r11,r14 ; x
```

Utilizamos el registro **rcx** como el *y₂-actual* (destino) y el registro **rdx** como el *x₂-actual* (destino). Recorremos la imagen destino de abajo hacia arriba, de izquierda a derecha. Al registro rcx debemos decrementarlo en uno porque arranca inicializado en cero.

```
mov rcx,r13 ; y2 = tamy
dec rcx      ; y2 = tamy-1
mov rdx,0    ; x2 = 0
```

En **r13** almacenamos el límite para r10, es decir tiene que recorrer el ciclo de y hasta que alcance su límite en y (offsety+tamy).

```
add r13,r15 ; r13 = offsety+tamy
```

En **rbx** almacenamos el límite para r11, es decir tiene que recorrer el ciclo de x hasta que alcance su límite en x (offsetx+tamx)

```
mov rbx,r12 ; rbx = tamx
add rbx,r14 ; rbx = offsetx+tamx
```

Ahora comienza el ciclo de iteración sobre la variable *y*. Cada vez que se ejecuta se comprueba que el *y-actual* (**r10**) sea menor que su límite (**r13**). Y por cada iteración se reinician los valores de *x* e *x₂* a la primer columna que debemos trabajar (Para la imagen fuente es el offset inicial de la variable *x*, para la imagen destino es el 0).

```
.loop_y:
cmp r10,r13      ; y < offsety+tamy
jge .endloop_y
mov r11,r14      ; x = offsetx
mov rdx,0        ; x2 = 0
```

Por consiguiente, para cada valor que vaya tomando la variable, se debe ejecutar un ciclo para poder iterar sobre todas las columnas (Ciclo de *x*). Por cada iteración, se compara si el *x-actual* (**r11**) es menor

que su límite (**rbx**). Si no lo es, se salta fuera del ciclo de x.

Luego, comienza a ejecutar el código propio del ciclo. Es necesario tener en cuenta que:

- ★ En **r8** está cargado el *row_size* de la imagen fuente.
- ★ En **r9** está cargado el *row_size* de la imagen de destino.

Vamos a calcular en cada iteración la cantidad de posiciones en memoria (bytes) que se le deben sumar a la posición de origen de la imagen para encontrarnos en la posición actual. Esto se almacena en el registro **rax**. Primero copiamos el contenido de **r10** (*y-actual*), lo multiplicamos por el largo de cada fila para así posicionarnos en la fila actual y por último le sumamos **r11** (que indica el número de fila actual) multiplicado por 4 porque cada Pixel tiene 4 bytes. De este modo, **rax** contiene el offset que debemos sumarle a la posición en memoria de la imagen fuente para situarnos en la posición actual.

Como nos encontramos utilizando registros Xmm, vamos a trabajar con 4 pixels a la vez, de modo que entran en un solo registro. Para levantar de memoria 4 pixels, utilizamos la instrucción *movdqu*.

Por último, lo que tenemos que hacer es guardar esos 4 pixels levantados con el mismo orden en la posición correspondiente de la imagen destino. Es decir, la misma columna *x* pero sin su offset (por eso mismo se usan dos registros distintos, ya que el offset de la imagen de destino es 0) y la fila y_2 va a ser la diferencia entre la cantidad de filas e *y*.

Análogamente a lo anterior, en **rax** se guarda el offset de la imagen destino, para acceder a memoria sumándoselo a la posición donde esta almacenada la imagen.

Se suma 4 a las variables **r11**(*x-actual*) y **rdx**(*x₂-actual*), porque en este paso avanzamos 4 pixels. Y el jump se ejecuta siempre, ya que la comparación se produce al principio del ciclo.

```
.loop_x:

cmp r11,rbx      ; x < offsetx+tamx
jge .endloop_x

mov rax,r10
imul rax,r8
lea rax,[rax+r11*4]

movdqu xmm0,[rdi+rax]

mov rax,rcx
imul rax,r9
lea rax,[rax+rdx*4]

movdqu [rsi+rax],xmm0

add r11,4
add rdx,4
jmp .loop_x
```

Una vez copiada toda una fila, debemos avanzar hacia la siguiente. De este modo: incrementamos en uno **r10**(*y-actual*) y **rcx**(*y₂-actual*).

Se ejecuta siempre un jump hacia el comienzo del ciclo en *y*, porque la comparación de ver si ya copiamos todas las filas se ejecuta al principio.

```
.endloop_x:
inc r10
dec rcx
jmp .loop_y
```

Una vez terminado el ciclo de *y*, sólo resta salir de la ejecución respetando la convención C.

```
.endloop_y:
Alinear la pila y popear todos los registros.
ret
```

3.2. Filtro Sierpinsky

Primero hacemos unos defines que vamos a utilizar luego.

```
offset: dd 3.0, 2.0, 1.0, 0.0
dq255:  dd 255.0,255.0,255.0,255.0
```

Limpiamos los registros xmm15 y los que vamos a utilizar como x e y.

```
Salvamos el tope de la pila y salvamos los registros a usar.
xor r15,r15 ; x
xor r14,r14 ; y
pxor xmm15,xmm15
```

Primero recorremos la imagen por filas (y) y luego por columnas (x). Para ello, primero comenzamos con el *ciclo de y*. En **ecx** estaba almacenada la cantidad total de filas de la imagen pasada por parámetro. Cuando alcanza este valor, finaliza el ciclo.

```
.loop_y:
    cmp r14d, ecx
    je .endloop_y    ; y = filas
```

Para entrar a este ciclo, x no debe haber alcanzado su tamaño máximo, es decir se repite hasta que hayamos trabajado con todas las columnas de la fila actual. Como para este filtro son necesarios los valores de x e y para hacer cálculos del coeficiente para cada iteración, nos debemos encargar de copiarlos en dos registros xmm transformados en puntos flotantes. Como vamos a trabajar en simultáneo, mediante la instrucción *shufps* y dándole como parámetro el 0h replicamos el valor 4 veces en el registro. Hacemos lo mismo tanto para y como para x, salvando que en x le debemos sumar el offset ya que los cuatro píxeles con los que vamos a trabajar en simultáneo pertenecen a la misma fila y a columnas contiguas.

```
.loop_x:
    cmp r15d,edx    ; x = cols
    je .endloop_x

    xor rax, rax
    mov rax, r15    ; rax = x
    cvtsi2ss xmm1,rax    ; xmm1 = (float) x
    shufps xmm1, xmm1, 0h ; replicó 4 veces rax en xmm1
    movups xmm2,[offset] ;
    addps xmm1, xmm2    ; le sumo el offset de x en los pixels (3,2,1,0)
    mov rax, r14    ; rax = y
    cvtsi2ss xmm2, rax ; xmm2 = (float) y
    shufps xmm2, xmm2, 0h ; replicó 4 veces rax en xmm2
```

Ahora, debemos multiplicar los valores de x e y por 255 y luego dividirlos por la cantidad de columnas y de filas respectivamente. Para hacer esto, mantenemos el empaquetamiento de cuatro valores por registro xmm. Luego de este paso, en **xmm1** va a quedar el valor del x actual, seguido de los valores de los tres siguientes, y en **xmm2** el valor del y actual replicado 4 veces.

```
movups xmm0,[dq255]
mulps xmm1,xmm0    ; xmm1 = x*255
mulps xmm2,xmm0    ; xmm2 = y*255

cvtsi2ss xmm3, rdx    ; xmm3 = cols
shufps xmm3, xmm3, 0h ; replicó 4 veces cols en xmm3

cvtsi2ss xmm4, rcx    ; xmm4 = filas
shufps xmm4, xmm4, 0h ; replicó 4 veces filas en xmm4

divps xmm1, xmm3    ; xmm1 = x*255/cols
divps xmm2, xmm4    ; xmm2 = y*255/filas
```

Lo que resta para calcular el coeficiente es hacer el xor y luego volver a dividir por 255. Para ejecutar el *pxor* es necesario volver a transformar los cuatro valores por registro en int, lo hacemos truncando. Luego se vuelve a transformar en puntos flotante de precisión simple y así dividirlos por 255.

```
cvttq2d xmm1,xmm1 ; convierte a int truncando
cvttq2d xmm2,xmm2

pxor xmm1,xmm2      ; Hace el xor entre los 4 Is y los 4 Js empaquetados

movups xmm0,[dq255]
cvtq2ps xmm1,xmm1   ; se convierte a float otra vez
divps xmm1,xmm0      ; se divide por 255.0
```

Como resultado, en **xmm1** tenemos los cuatro coeficientes que corresponden a los próximos cuatro píxeles que vamos a procesar. Hacemos una copia de los coeficientes al registro **xmm15**, porque el registro **xmm1** lo vamos a utilizar para guardar las lecturas en memoria. En **rax** hacemos el cálculo del offset en la matriz imagen que corresponde a los *x* e *y* actuales. Leemos los valores de *x* e *y* actuales, multiplicamos al *y* por el *row_size* de la imagen fuente y al valor de *x* lo multiplicamos por cuatro porque por cada píxel son cuatro bytes. En **rdi** se encuentra el puntero al origen de la imagen, entonces para hacer la lectura simplemente sumamos *rax* con *rdi* y lo guardamos en **xmm0**, haciendo una copia en **xmm1**.

```
movups xmm15,xmm1
; float s = (i xor j) / 255.0;

xor rax,rax
mov r12d,r14d ; r12 = y
mov r13d,r15d ; r13 = x
imul r12d,r9d ; r12 = y * row_size_src
imul r13d,4   ; r13 = x * 4
mov eax,r12d ;
add eax,r13d ; rax = y*row_size_src + x*4

movdqu xmm0,[rdi+rax] ; Lectura de 16 Bytes del source
movdqu xmm1,xmm0      ; Se hace una copia
```

Una vez que ya tenemos copiado en dos registros los píxeles con los que vamos a trabajar, necesitamos convertirlos en puntos flotantes. Para ello es necesario primero expandirlos y por este motivo se desempaquetan mediante las instrucciones *punpcklbw* y *punpcklwd* con el registro **xmm2** que fue limpiado. De este modo es posible transformarlos en puntos flotantes mediante la instrucción *cvtq2ps*, ya que tenemos un píxel por registro y esto es 4 números en punto flotante (4 Bytes).

```
pxor xmm2,xmm2

punpcklbw xmm0,xmm2 ; parte baja
movdqu xmm6,xmm0    ; copia de la parte baja
punpcklwd xmm0,xmm2 ; baja de la baja 0
punpckhbw xmm6,xmm2 ; alta de la baja 1

punpckhbw xmm1,xmm2 ; parte alta
movdqu xmm7,xmm1    ; copia de la parte alta
punpckhwd xmm1,xmm2 ; alta de alta 3
punpcklwd xmm7,xmm2 ; baja de la alta 2

cvtq2ps xmm0,xmm0
cvtq2ps xmm6,xmm6
cvtq2ps xmm1,xmm1
cvtq2ps xmm7,xmm7
```

En **xmm15** habíamos copiado los coeficientes, ahora que ya tenemos los cuatro píxeles con sus valores en puntos flotantes sólo resta multiplicarlos. Primero hacemos una copia de **xmm15** en **xmm14** porque lo vamos a pisar, ya que en cada caso necesitamos replicar cuatro veces los cuatro valores de coeficientes almacenados. De este modo, los tres canales de cada píxel quedan multiplicados con su mismo coeficiente. Mediante *shufps* replicamos el valor que necesitamos, tal cual figura comentado en binario.

```
movdqu xmm14,xmm15
shufps xmm15,xmm15, 0x00 ; 00 00 00 00
mulps xmm1,xmm15

movdqu xmm15,xmm14
shufps xmm15,xmm15, 0x55 ; 01 01 01 01
mulps xmm7,xmm15

movdqu xmm15,xmm14
shufps xmm15,xmm15, 0xAA ; 10 10 10 10
mulps xmm6,xmm15

movdqu xmm15,xmm14
shufps xmm15,xmm15, 0xFF ; 11 11 11 11
mulps xmm0,xmm15
```

Una vez ya calculados todos los valores finales, debemos transformarlos en enteros y empaquetarlos al tamaño de char como era originalmente para poder asignárselos a la imagen destino. Mediante la instrucción *cvttps2dq* se convierten los cuatro valores de cada píxel, que estaban en puntos flotantes de precisión simple a enteros empaquetados. Luego, mediante la instrucción *packusdw* empaquetamos de a dos registros, para así modificar le tamaño de ints a shorts. Por último, unimos los cuatro píxeles en un registro con la instrucción *packusdw* que los transforma de shorts a char. Todo esto lo realizamos mediante unsigned saturation para eliminar casos de exceso que puedan hacer que la fotografía final no quede como lo deseado. Debemos volver a calcular **rax** para la imagen de destino, ya que no necesariamente tengan el mismo *row_size*, el código es el mismo que el utilizado anteriormente, sólo que en **r8d** se encuentra el tamaño de *row_size_destino*. En **rsi** se encuentra el puntero a la imagen de destino, por este motivo se suman ambos registros y se asigna a esa posición de memoria lo calculado que se encuentra en el registro **xmm0**. Avanzamos 4 en *x*, porque leímos cuatro píxeles y saltamos al comienzo del ciclo en *x* otra vez.

```
cvttps2dq xmm1,xmm1
cvttps2dq xmm0,xmm0
cvttps2dq xmm6,xmm6
cvttps2dq xmm7,xmm7

packusdw xmm0,xmm6
packusdw xmm7,xmm1

packuswb xmm0,xmm7

xor rax,rax
mov r12d,r14d ; r12 = y
mov r13d,r15d ; r13 = x
imul r12d,r8d ; y * row_size_dest
imul r13d,4 ; x * 4
mov eax,r12d ;
add eax,r13d ; rax = y*row_size_dest + x*4

movdqa [rsi+rax],xmm0

add r15,4 ; avanza 4 sobre x
jmp .loop_x
```


Para culminar el ciclo de x , reiniciamos a cero el valor de x y aumentamos en uno el valor de y , volviendo al comienzo del ciclo de y .

```
.endloop_x:
    inc r14      ; y++
    xor r15,r15 ; reinicio x
    jmp .loop_y
```

Cuando terminó el ciclo de y , terminó la ejecución del filtro. Por lo tanto, popeamos los registros usados y retornamos.

```
.endloop_y:
    Hacemos pop de los registros salvados
    ret
```

3.3. Filtro Bandas

Algunos defines que luego vamos a necesitar

```
color1 dw 64,64,64,64,64,64,64,64
color3 db 128,128,128,255
color4 db 192,192,192,255
color5 db 255,255,255,255

all_64w dw -64,-64,-64,-64,-64,-64,-64,-64

color1_bound dw 96,96,96,96,96,96,96,96
color2_bound dw 288,288,288,288,288,288,288,288
color3_bound dw 480,480,480,480,480,480,480,480
color4_bound dw 672,672,672,672,672,672,672,672

rgb_only_mask dd 0x00FFFFFF,0x00FFFFFF,0x00FFFFFF,0x00FFFFFF

all_1_mask dd 0xFFFFFFFF,0xFFFFFFFF,0xFFFFFFFF,0xFFFFFFFF
all_0_mask dd 0x00000000,0x00000000,0x00000000,0x00000000

magic_shuffle db 8,8,8,8,10,10,10,10,0,0,0,0,2,2,2,2
```

Salvamos los registros a utilizar, para respetar la convención C. Luego limpiamos los registros que vamos a usar como *x* e *y* y el registro *xmm15*.

```
Pusheamos los registros que vamos a utilizar
xor r15,r15 ; x
xor r14,r14 ; y
pxor xmm15,xmm15
```

Vamos a recorrer la imagen primero por filas y luego por columnas, por este motivo comenzamos con el ciclo de *y*, el cual termina cuando el *y* actual (**r14**) alcance su valor máximo de filas que está almacenado en **ecx**.

```
.loop_y:
    cmp r14d,ecx
    je .endloop_y ; y = filas
```

Comienza el ciclo de *x*, el cual termina cuando **r15** alcanza el valor máximo de columnas el cual está almacenado en **edx**. Luego, en **rax** vamos a calcular el offset dentro de la matriz imagen sumando el valor de *y* actual multiplicado por el *row_size* fuente con el valor de *x* actual multiplicado por cuatro porque por cada píxel hay cuatro bytes.

```
.loop_x:
    cmp r15d,edx
    je .endloop_x ; x = cols

    ; armo el offset
    xor rax,rax
    mov r12d,r14d ; r12 = y
    mov r13d,r15d ; r13 = x
    imul r12d,r9d ; y * row_size_src
    imul r13d,4 ; x * 4
    mov eax,r12d ;
    add eax,r13d ; rax = y*row_size_src + x*4
```

Luego de haber calculado el offset actual, leemos de memoria 16 Bytes y los almacenamos en **xmm1**, copiándolos también en **xmm2**. Para evitar sumar número que no correspondan, ponemos en cero al canal

alpha mediante el *pand* y la etiqueta definida antes *rgb_only_mask*. De la lectura que hicimos, guardamos en **xmm1** la parte baja y en **xmm2** la parte alta desempaquetando. Luego de eso, debemos ejecutar la suma horizontal mediante la instrucción *phaddw*, así podemos sumar los cuatro bits de cada píxel. Para esto debemos ejecutarlas dos veces, lo hacemos en distinto orden así luego quedan intercalados entre los registros **xmm1** y **xmm2** y podemos trabajar con los 4 píxeles a la vez si los sumamos.

```
movdqu xmm1,[rdi+rax] ; agarro 16 bytes del source
movdqu xmm2,xmm1

movdqu xmm15,[rgb_only_mask]
pand xmm1,xmm15
pand xmm2,xmm15

pxor xmm7,xmm7
punpcklbw xmm1,xmm7 ; baja 0
punpckhbw xmm2,xmm7 ; alta 8

pxor xmm7,xmm7
pxor xmm8,xmm8

phaddw xmm1,xmm7
phaddw xmm7,xmm1

movdqu xmm1,xmm7

phaddw xmm2,xmm8
phaddw xmm2,xmm8

paddw xmm1,xmm2
```

Copiamos en los registros **xmm14**, **xmm13** y **xmm12** las etiquetas predefinidas. Luego en **xmm15** vamos a ir copiando la etiqueta que necesitemos en cada uno de los cuatro casos. El comportamiento en cada uno de los cuatro casos va a ser similar:

- ▷ Hacemos una copia de la suma original obtenida.
- ▷ En el registro **xmm15** guardamos el valor máximo para cada banda.
- ▷ Comparamos para cada uno de los cuatro valores mediante la instrucción *pcmpgtw*. Si pertenece a la banda actual, va a poner en uno el valor del número.
- ▷ Los valores que hayan dado *true* le ponemos todos sus bits en 1, así después podemos hacer un *pand*.
- ▷ Hacemos el and con el *color_bound* así quedan los píxeles de la banda del color que deban ser y el resto en cero.
- ▷ Acumulamos lo obtenido en distintos registros.

Al finalizar estas cuatro ejecuciones, sumamos los cuatro acumuladores con el fin de tener la versión final.

```
movdqu xmm14,[all_64w]
movdqu xmm13,[all_1_mask]
movdqu xmm12,[color1]

movdqu xmm0,xmm1 ; salvo una copia
movdqu xmm15,[color1_bound]
pcmpgtw xmm15,xmm1 ; greater than.
pandn xmm15,xmm13 ; deja todos 1 en los lugares que va el color_1
pand xmm15,xmm12

movdqu xmm1,xmm15 ; acumulo lo que tengo en xmm1

movdqu xmm2,xmm0 ; salvo una copia
movdqu xmm15,[color2_bound]
pcmpgtw xmm15,xmm2
pandn xmm15,xmm13
pand xmm15,xmm12

movdqu xmm2,xmm15

movdqu xmm3,xmm0
movdqu xmm15,[color3_bound]
pcmpgtw xmm15,xmm3
pandn xmm15,xmm13
pand xmm15,xmm12

movdqu xmm3,xmm15

movdqu xmm4,xmm0
movdqu xmm15,[color4_bound]
pcmpgtw xmm15,xmm4
pandn xmm15,xmm13
pand xmm15,xmm12

movdqu xmm4,xmm15

paddusb xmm1,xmm2
paddusb xmm1,xmm3
paddusb xmm1,xmm4
```

Como para agilizar el comportamiento debimos mezclar los píxeles, debemos reordenarlos como estaban originalmente. En este momento usamos la etiqueta *magic_shuffle* y mediante la instrucción *pshufb* quedan almacenados en el registro **xmm1**. Sólo resta calcular el offset de la imagen de destino en **rax**, lo hacemos de manera análoga al cálculo del offset.fuente, sólo que ahora tenemos en cuenta el **row_size_dest**. Guardamos lo calculado en la memoria de la foto del destino, aumentamos 4 el contador de **x** y volvemos al ciclo de **x**.

```
; reordeno los bytes a como estaba original
movdqu xmm15,[magic_shuffle]
pshufb xmm1,xmm15

; armo el offset
xor rax,rax
mov r12d,r14d ; r12 = y
mov r13d,r15d ; r13 = x
imul r12d,r8d ; y * row_size_dest
imul r13d,4 ; x * 4
mov eax,r12d ;
add eax,r13d ; rax = y*row_size_dest + x*4

movdqu [rsi+rax],xmm1

add r15,4 ; avanzo 4 sobre x
jmp .loop_x
```

El ciclo de x finaliza cuando x ya alcanzó el valor máximo de columnas. Por esto, incrementamos en uno la fila y reiniciamos x a cero, volviendo a empezar el ciclo de y .

```
.endloop_x:
    inc r14
    xor r15,r15 ; reinicio x
    jmp .loop_y
```

Al finalizar el ciclo de y , no resta hacer nada. Por lo tanto, para respetar la convención C, hacemos pop de los registros usados y retornas con *ret*.

```
.endloop_y:
    Pop de los registros usados
    ret
```

3.4. Filtro Motion Blur

Primero hacemos algunos define que vamos a necesitar luego:

```
cerocomados    dd    0.2 , 0.2 , 0.2 , 0.2
rgb_only       dd    0x00FFFFFF, 0x00FFFFFF, 0x00FFFFFF, 0x00FFFFFF
```

Salvamos la base de la pila y los registros a utilizar. Alineamos la pila.
Limpiamos con xor o pxor los registros a utilizar: r10, r11, r15, r13 y xmm13.

```
; r10 = i
; r11 = j
mov r15d,edx ; r15 = cols
mov r14d,ecx ; r14 = filas
sub r15,2    ; r15 = cols -2
sub r14,2    ; r14 = filas - 2
```

El registro **xmm13** es limpiado para luego utilizarlo a la hora de asignarle el borde negro. Se les resta 2 a las filas y a las columnas porque los últimos dos píxeles del borde van de negro.

Comienza el ciclo de y, es decir, primero recorremos por filas y por cada fila recorremos por columnas. Se recorren todas las filas y cada vez que se entra en el *loop_x* se inicializa la fila en cero.

```
.loop_y:
    cmp r10d,ecx ; i < filas
    jge .endloop_y

    mov r11,0    ; j = 0
```

El *loop_x* se recorre hasta que j alcanza el valor de la última columna. En el registro **rax** vamos a asignarle el número de orden del byte a tratar, es decir la fila actual multiplicada por la cantidad de columnas sumado a la columna actual multiplicada por cuatro porque es lo que ocupa un píxel. Antes de asignarle el nuevo valor a la imagen, debemos comprobar si el píxel a tratar es de borde en cuyo caso queda negro (.cero) o si no lo es (.nocero). Todos los píxeles que son bordes son los que están ubicados en las primeras dos filas o columnas y los que están ubicados en las últimas dos filas o columnas.

```
.loop_x:
    cmp r11d,edx ; j < cols
    jge .endloop_x

    mov rax,r10  ; rax = i
    imul eax,r8d ; rax = i * row_size
    lea rax,[rax+r11*4] ; rax = i * row_size + j * 4

    cmp r11,2
    jl .cero     ; j < 2
    cmp r10,2
    jl .cero     ; i < 2
    cmp r10,r14
    jge .cero    ; i >= cols -2
    cmp r11,r15
    jge .cero    ; j >= filas -2

    jmp .nocero
```

Si el píxel a tratar debe quedar en negro (cero) lo que debemos hacer es asignarle ceros mediante el registro **xmm13**:

```
.cero:
    movq [rsi+rax],xmm13
    add r11,2
    jmp .loop_x
```

A continuación comienza la rutina a llevar a cabo en caso de que el píxel a tratar no sea borde. Se copia en **rbx** el puntero actual con el que vamos a trabajar. Recordar que en **rax** teníamos el offset en la imagen y en **rdi** el puntero al comienzo de la misma. Vamos a ejecutar 4 píxeles por vez (porque son los que entran en un registro xmm) y para ello necesitamos hacer cinco lecturas a memoria y así tener para cada píxel sus cuatro vecinos que van a influir en su valor final. En **xmm14** copiamos el contenido de **rgb_only** el cual nos permitirá sólo trabajar con los canales RGB de cada píxel más adelante.

```
.nocero:

    lea rbx,[rdi+rax]
    movdqu xmm1,[rbx]          ; (i,j)
    movdqu xmm2,[rbx+r8*1+4]   ; (i+1,j+1)
    movdqu xmm3,[rbx+r8*2+8]   ; (i+2,j+2)
    mov r12,rbx
    sub r12,r8
    movdqu xmm4,[r12-4]        ; (i-1,j-1)
    sub r12,r8
    movdqu xmm5,[r12-8]        ; (i-2,j-1)

    movdqu xmm14,[rgb_only]
```

Para continuar necesitamos desempaquetar los valores para poder hacer las cuentas, ya que necesitamos los valores de los píxeles agrupados en números de mayor tamaño para el cual contemos con sus instrucciones necesarias y también un orden de precisión mayor. Para ello, los expandimos con ceros que van a provenir de haber limpiado el registro **xmm0**.

```
pxor xmm0,xmm0
movdqu xmm6,xmm1          ; xmm6 = (i,j) || (i,j+1) || (i,j+2) || (i,j+3)
punpcklbw xmm1,xmm0       ; xmm1 = (i,j) || (i,j+1)
punpckhbw xmm6,xmm0       ; xmm6 = (i,j+2) || (i,j+3)

movdqu xmm7,xmm2          ; xmm7 = (i+1,j+1) || (i+1,j+2) || (i+1,j+3) || (i+1,j+4)
punpcklbw xmm2,xmm0       ; xmm2 = (i+1,j+1) || (i+1,j+2)
punpckhbw xmm7,xmm0       ; xmm7 = (i+1,j+3) || (i+1,j+4)

movdqu xmm8,xmm3          ; xmm8 = (i+2,j+2) || (i+2,j+3) || (i+2,j+4) || (i+2,j+5)
punpcklbw xmm3,xmm0       ; xmm3 = (i+2,j+2) || (i+2,j+3)
punpckhbw xmm8,xmm0       ; xmm8 = (i+2,j+4) || (i+2,j+5)

movdqu xmm9,xmm4          ; xmm9 = (i-1,j-1) || (i-1,j) || (i-1,j+1) || (i-1,j+2)
punpcklbw xmm4,xmm0       ; xmm4 = (i-1,j-1) || (i-1,j)
punpckhbw xmm9,xmm0       ; xmm9 = (i-1,j+1) || (i-1,j+2)

movdqu xmm10,xmm5         ; xmm10 = (i-2,j-1) || (i-2,j) || (i-2,j+1) || (i-2,j+2)
punpcklbw xmm5,xmm0       ; xmm5 = (i-2,j-1) || (i-2,j)
punpckhbw xmm10,xmm0      ; xmm10 = (i-2,j+1) || (i-2,j+2)
```

Por cada registro tenemos dos píxeles, de este modo ya podemos realizar la suma entre los cinco píxeles vecinos que luego de dividirlo por la constante van a ser asignados al píxel central de estos cinco. De este modo, contamos con la instrucción *paddw* la cual suma empaquetadamente de a Word (2 bytes). Utilizamos los registros **xmm1** y **xmm6** como acumuladores de la suma.

```
paddw xmm1,xmm2      ; xmm1 = (i,j)+(i+1,j+1) || (i,j+1)+(i+1,j+2)
paddw xmm1,xmm3      ; xmm1 = (i,j)+(i+1,j+1)+(i+2,j+2) || (i,j+1)+(i+1,j+2)+(i+2,j+3)
paddw xmm1,xmm4      ; xmm1 = (i,j)+(i+1,j+1)+(i+2,j+2)+(i-1,j-1) ||
                      ; (i,j+1)+(i+1,j+2)+(i+2,j+3)+(i-1,j)
paddw xmm1,xmm5      ; xmm1 = (i,j)+(i+1,j+1)+(i+2,j+2)+(i-1,j-1)+(i-2,j-1) ||
                      ; (i,j+1)+(i+1,j+2)+(i+2,j+3)+(i-1,j)+(i-2,j)

paddw xmm6,xmm7      ; xmm6 = (i,j+2)+(i+1,j+3) || (i,j+3)+(i+1,j+4)
paddw xmm6,xmm8      ; xmm6 = (i,j+2)+(i+1,j+3)+(i+2,j+4) || (i,j+3)+(i+1,j+4)+(i+2,j+5)
paddw xmm6,xmm9      ; xmm6 = (i,j+2)+(i+1,j+3)+(i+2,j+4)+(i-1,j+1) ||
                      ; (i,j+3)+(i+1,j+4)+(i+2,j+5)+(i-1,j+2)
paddw xmm6,xmm10     ; xmm6 = (i,j+2)+(i+1,j+3)+(i+2,j+4)+(i-1,j+1)+(i-2,j+1) ||
                      ; (i,j+3)+(i+1,j+4)+(i+2,j+5)+(i-1,j+2)+(i-2,j+2)
```

Para poder multiplicarlos por una constante debemos transformar los valores en puntos flotantes, para ello primero debemos desempaquetarlos ampliandolos con cero (adquiriendo así un tamaño de Doubleword) y luego mediante la instrucción *cvt dq2ps* convertirlos de enteros con signo a puntos flotantes de precisión simple sin perder el empaquetamiento que nos permite trabajar con los tres valores del píxel a la vez. (En realidad son cuatro, pero un canal no lo usamos).

```
movdqu xmm2,xmm1      ; xmm2 = (i,j)+(i+1,j+1)+(i+2,j+2)+(i-1,j-1)+(i-2,j-1) ||
                      ; (i,j+1)+(i+1,j+2)+(i+2,j+3)+(i-1,j)+(i-2,j)
movdqu xmm7,xmm6      ; xmm7 = (i,j+2)+(i+1,j+3)+(i+2,j+4)+(i-1,j+1)+(i-2,j+1) ||
                      ; (i,j+3)+(i+1,j+4)+(i+2,j+5)+(i-1,j+2)+(i-2,j+2)

punpcklwd xmm1,xmm0    ; xmm1 = (i,j)+(i+1,j+1)+(i+2,j+2)+(i-1,j-1)+(i-2,j-1)
punpckhwd xmm2,xmm0    ; xmm2 = (i,j+1)+(i+1,j+2)+(i+2,j+3)+(i-1,j)+(i-2,j)

punpcklwd xmm6,xmm0    ; xmm6 = (i,j+2)+(i+1,j+3)+(i+2,j+4)+(i-1,j+1)+(i-2,j+1)
punpckhwd xmm7,xmm0    ; xmm7 = (i,j+3)+(i+1,j+4)+(i+2,j+5)+(i-1,j+2)+(i-2,j+2)

cvt dq2ps xmm1,xmm1
cvt dq2ps xmm2,xmm2
cvt dq2ps xmm6,xmm6
cvt dq2ps xmm7,xmm7
```

Ahora que los valores ya están transformados en puntos flotantes, podemos hacer la multiplicación (también empaquetada). De este modo, calculamos todo un píxel por cada registro. Una vez hechas las multiplicaciones, volvemos a convertir los valores en enteros mediante la instrucción *cvtps2dq*. Como habíamos definido antes, *cerocomados* guarda el valor 0.2 cuatro veces.

```
movdqu xmm15,[cerocomados]

mulps xmm1,xmm15
mulps xmm2,xmm15
mulps xmm6,xmm15
mulps xmm7,xmm15

cvtps2dq xmm1,xmm1
cvtps2dq xmm2,xmm2
cvtps2dq xmm6,xmm6
cvtps2dq xmm7,xmm7
```


Ahora resta volver a empaquetar para que queden los cuatro píxeles, que estábamos trabajando en simultáneo, en un sólo registro. Primero juntamos el valor final de **xmm1** (i, j) con el de **xmm2** ($i, j + 1$) y el de **xmm6** ($i, j + 2$) con el del **xmm7** ($i, j + 3$) y por último unimos todos en un mismo registro (**xmm1**). Guardamos el contenido de xmm1 en memoria, a la dirección calculada previamente en **rax**, sumada al origen del puntero destino (**rsi**). Sumamos cuatro al contador de **r11** (columnas) porque por cada lectura en memoria, trabajamos con 4 píxeles, y luego volvemos al ciclo de x.

```
packusdw xmm1,xmm2
packusdw xmm6,xmm7

packuswb xmm1,xmm6

movdqu [rsi+rax],xmm1

add r11,4
jmp .loop_x
```

Al terminar el ciclo de x, lo que hacemos es incrementar en uno **r10**, el cual mide la fila actual. Es decir, avanzamos de fila y saltamos al ciclo de y.

```
.endloop_x:
    inc r10
    jmp .loop_y
```

Al terminar el ciclo de y, queda finalizada la ejecución del filtro. Por lo tanto, solo queda restaurar los registros y retornar.

```
.endloop_y:
    Hacemos pop de los registros salvados
    ret
```

4. Enunciado y solución

4.1. Mediciones

Realizar una medición de performance *rigurosa* es más difícil de lo que parece. En este experimento deberá realizar distintas mediciones de performance para verificar que sean buenas mediciones.

En un sistema “ideal” el proceso medido corre solo, sin ninguna interferencia de agentes externos. Sin embargo, una PC no es un sistema ideal. Nuestro proceso corre junto con decenas de otros, tanto de usuarios como del sistema operativo que compiten por el uso de la CPU. Esto implica que al realizar mediciones aparezcan “ruidos” o “interferencias” que distorsionen los resultados.

El primer paso para tener una idea de si la medición es buena o no, es tomar varias muestras. Es decir, repetir la misma medición varias veces. Luego de eso, es conveniente descartar los outliers¹, que son los valores que más se alejan del promedio. Con los valores de las mediciones resultantes se puede calcular el promedio y también la varianza, que es algo similar al promedio de las distancias al promedio².

Las fórmulas para calcular el promedio μ y la varianza σ^2 son

$$\mu = \frac{1}{n} \sum_{i=1}^n x_i \quad \sigma^2 = \frac{\sum_{i=1}^n (x_i - \mu)^2}{n}$$

4.2. Filtro *cropflip*

Programar el filtro *cropflip* en lenguaje C y luego en ASM haciendo uso de las instrucciones vectoriales (SSE).

Experimento 1.1 - análisis el código generado

En este experimento vamos a utilizar la herramienta `objdump` para verificar como el compilador de C deja ensamblado el código C.

Ejecutar

```
objdump -Mintel -D cropflip_c.o
```

¿Cómo es el código generado? Indicar

a) Por qué cree que hay otras funciones además de `cropflip_c` b) Cómo se manipulan las variables locales c) Si le parece que ese código generado podría optimizarse

Muchas de las funciones generadas se agregan al compilar usando parámetros de debugging.

El cambio que puede notarse a simple vista es que el tamaño del código generado con optimización es considerablemente menor al generado sin optimizaciones, un 34 % menor. Por otro lado, calcula una única vez la posición de la cual deberá leer el pixel a procesar, ya que es una operación que se realiza cuatro veces por cada lectura y escritura de un pixel.

El código optimizado también disminuye el uso del stack para variables locales y en su lugar utiliza registros, por esa razón necesita armar el stackframe y pushear los registros que la convención C pide preservar. Esto no se da en el código sin optimizar, todas las variables locales se guardan en el stack, se asignan a algún registro libre al momento de usarla y luego vuelve a ser guardada en el stack. Si bien la cantidad de accesos a memoria es sustancialmente mayor en el código sin optimizar por el hecho de que se accede constantemente a las variables locales alojadas en el stack, también hay que considerar que al haber tantos accesos al stack el hit-rate de la caché debería ser bastante alto, agilizando la lectura, no así la escritura.

Experimento 1.2 - optimizaciones del compilador

Compile el código de C con flags de optimización. Por ejemplo, pasando el flag `-O1`³. Indicar

¹en español, valor atípico: http://es.wikipedia.org/wiki/Valor_atpico

²en realidad, elevadas al cuadrado en vez de tomar el módulo

³agregando este flag a `CCFLAGS64` en el `makefile`

1. Qué optimizaciones observa que realizó el compilador
2. Qué otros flags de optimización brinda el compilador
3. Los nombres de tres optimizaciones que realizan los compiladores.

GCC provee un arsenal de optimizaciones disponibles para que podamos usar cuando lo creamos conveniente, asimismo provee unos flags para poder compilar el código con un conjunto de optimizaciones, estos son O1, O2, O3, Os, entre otros. O1, O2 y O3 son optimizaciones generales no agresivas, es decir, no deberían modificar el funcionamiento del programa. Otras optimizaciones podrían asumir que todas las operaciones aritméticas son sin signo o que están bien implementadas y no es necesario chequear la consistencia de los resultados (fno-math-errno). Estas optimizaciones tienen como objetivo mejorar la performance del programa, sin embargo Os tiene como objetivo reducir el tamaño del ejecutable, es decir, activa todas las optimizaciones que contribuyen a reducir la cantidad de instrucciones regardless del impacto que esto pueda tener en la performance del programa.

Concretamente al activar las optimizaciones O3 en *cropflip* el cambio más notorio es que GCC utiliza instrucciones SSE para procesar 4 pixels por iteración, lo cual le brinda un boost de velocidad impresionante. Además de haberse reducido la cantidad de instrucciones.

Experimento 1.3 - calidad de las mediciones

1. Medir el tiempo de ejecución de *cropflip* 10 veces. Calcular el promedio y la varianza. Consideraremos outliers a los 2 mayores tiempos de ejecución de la medición y también a los 2 menores, por lo que los descartaremos. Recalcular el promedio y la varianza después de hacer este descarte. Realizar un gráfico que presente estos dos últimos items.

Luego de ejecutar 10 veces el filtro *Cropflip* obtuvimos los siguientes resultados:

ASM	C
70.925	1.152.187
70.521	1.151.544
32.859	761.937
43.720	649.248
64.236	1.152.847
70.793	1.153.061
71.271	1.152.765
44.616	1.152.798
56.124	725.420
71.775	1.155.718
Esperanza	
59.684	1.020.752,5
Desvío estándar	
13.720,5329	203.626,443

El cuadro denota la cantidad de ciclos de clock utilizada por cada ejecución del programa.

Se puede apreciar que la esperanza (promedio de las mediciones) bajo Lenguaje C es notablemente mayor para las diez mediciones llevadas a cabo. Además poseen una desviación estándar también mayor. La desviación estándar mide el grado de dispersión o variabilidad, por lo que demuestra menor estabilidad en las mediciones bajo Lenguaje C. Por lo tanto podemos afirmar que bajo el contexto de nuestras mediciones el tiempo de ejecución del mismo programa en lenguaje ASM posee un tiempo de cómputo menor.

Luego de eliminar los dos valores más altos y los dos valores más bajos, recalculamos obteniendo los siguientes datos:

Esperanza: 62.869,16667(ASM) y 1.087.346,333(C)

Desvío estándar: 9.719,205547(ASM) y 145.528,2027(C)

Se puede ver que al eliminar los outliers, el valor de la esperanza aumenta en ambos casos. Al haber hecho una medición tan pequeña, los cálculos de la esperanza no son tan estables y no son los más “parecidos” a la realidad. Lo que sí ocurre es que al eliminar los valores mas alejados del promedio, como era de esperar, disminuye la desviación estándar al haber dejado sólo las mediciones más cercanas entre sí.

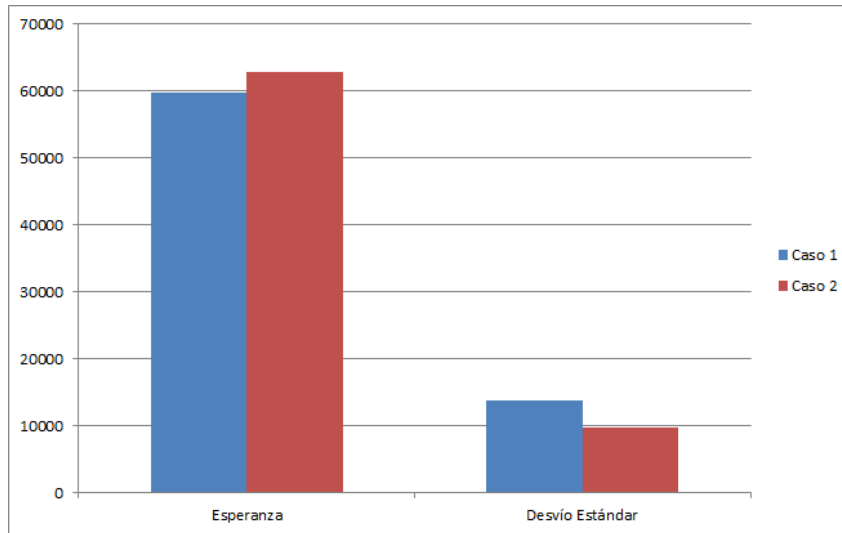


Figura 1: Assembler

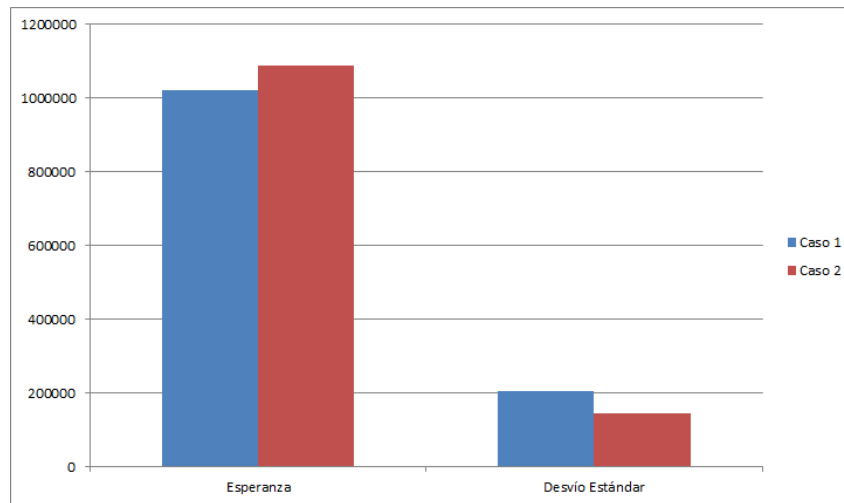


Figura 2: C

Siendo el Caso 1 las mediciones de esperanza y Desvío estándar para todos los casos de test y el Caso 2 las mediciones sin tener en cuenta los cuatro outliers.

- Implementar un programa en C que no haga más que ciclar infinitamente sumando 1 a una variable. Lanzar este programa tantas veces como *cores lógicos* tenga su procesador. Medir otras 10 veces mientras estos programas corren de fondo. Realizar los mismos casos de experimentación que en el ejercicio anterior.

Los resultados obtenidos en esta experimentación fueron menores que los anteriores:

ASM	C
33.585	542.928
33.798	544.155
33.402	544.857
33.228	543.687
33.159	543.252
33.441	543.324
34.089	544.224
33.768	760.359
34.563	542.448
34.473	542.982
Esperanza	
33.750,6	565.221,6
Desvío estándar	
465,49	65.049,34

Luego de eliminar los dos valores más altos y los dos valores más bajos, recalculamos obteniendo los siguientes datos:

Esperanza: 33.680,5 (ASM) y 543.604(C)

Desvío estándar: 235,364 (ASM) y 462,615(C)

Acá también se puede apreciar que al eliminar los outliers, el Desvío estándar disminuye su valor.

Se puede observar que las mediciones mejoran con la ejecución del ciclo infinito de fondo, esto se debe a que fue ejecutado en una computadora con un procesador i5. La familia ix de intel varía la frecuencia del clock de acuerdo al uso requerido. Por eso cuando ejecutamos los sumadores puede llegar a

tardar menos que en la versión sin ellos.

Por este motivo, volvimos a ejecutar este caso una mayor cantidad de veces para que el procesador no modifique la frecuencia de clock, obteniendo los siguientes resultados:

- Teniendo en cuenta todas las mediciones realizadas:
Esperanza: 434.338,4653(ASM) y 12.183.006,91(C).
Desvío estándar: 659,0435989(ASM) y 8.230.692,407(C).
- Eliminando los *outliers*:
Esperanza: 270.056,844(ASM) y 9.525.801,857(C).
Desvío estándar: 524,7013438(ASM) y 8.648.990,473(C).

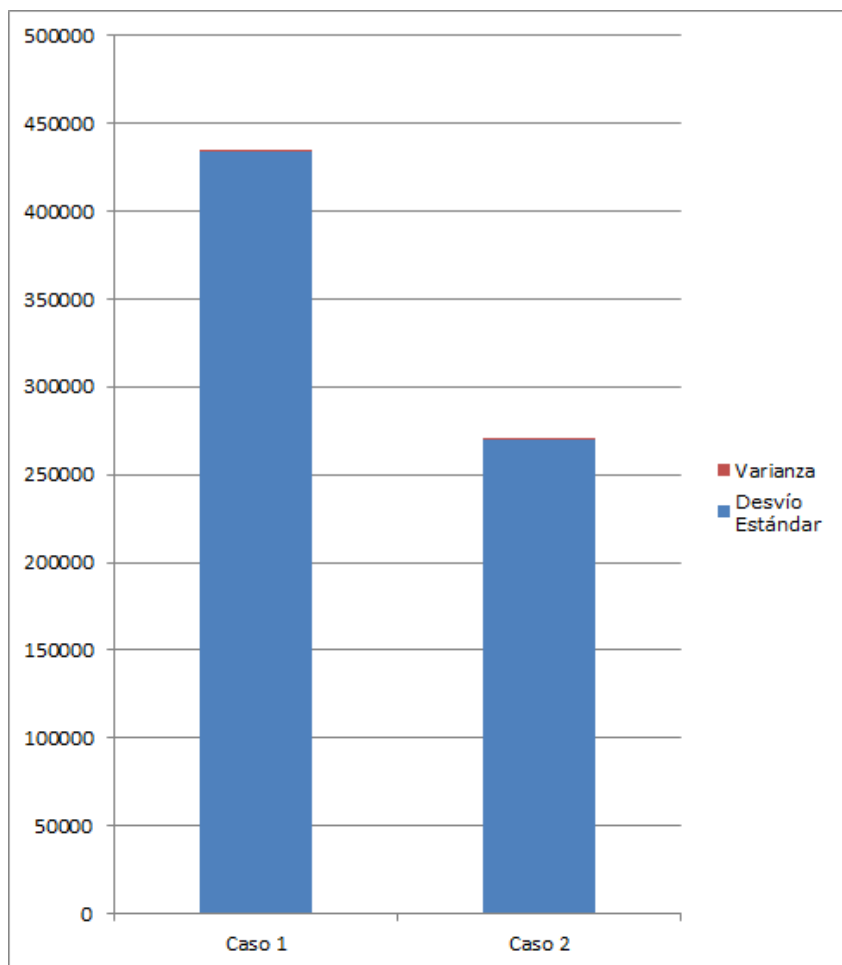


Figura 3: Assembler

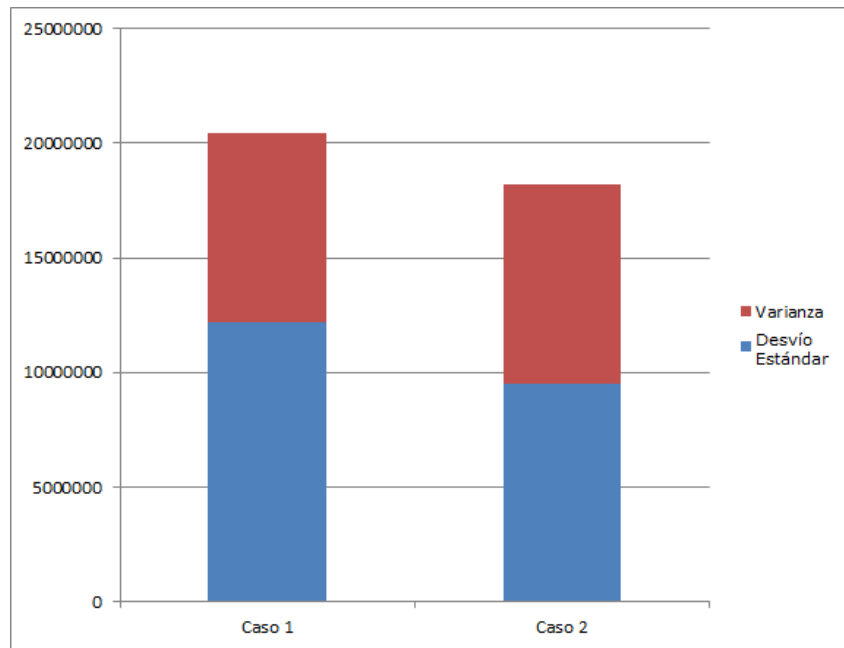


Figura 4: C

Siendo el Caso 1 las mediciones de esperanza y Desvío estándar para todos los casos de test y el Caso 2 las mediciones sin tener en cuenta los cuatro outliers.

Tal como se puede apreciar en los gráficos, la varianza en Assembler es casi despreciable pero la de C no lo es. Sin embargo, ahora sí ocurre algo más intuitivo que es que aumentaron los valores de la esperanza en proporciones notables.

A partir de aquí todos los experimentos de mediciones deberán hacerse igual que en el presente ejercicio: tomando 10 mediciones, luego descartando outliers y finalmente calculando promedio y Desvío estándar.

Decidimos:

Realizar 12000 mediciones por experimento, eliminando los primeros dos mil quinientos casos que hayan llevado menos ciclos de clock y los dos mil quinientos casos que hayan llevado la mayor cantidad de ciclos de clock.

Lo determinamos de esta manera, ya que dejar dos mediciones afuera, como dice el enunciado, no tiene influencia en los cálculos de la esperanza y la varianza para muestras tan grandes.

Luego de experimentar distintas cantidades de casos de testeo, notamos que elegir 7000 valores pertenecientes a la franja del medio de los 12000 es una solución lo suficientemente estable, por lo cual es la que llevamos a cabo.

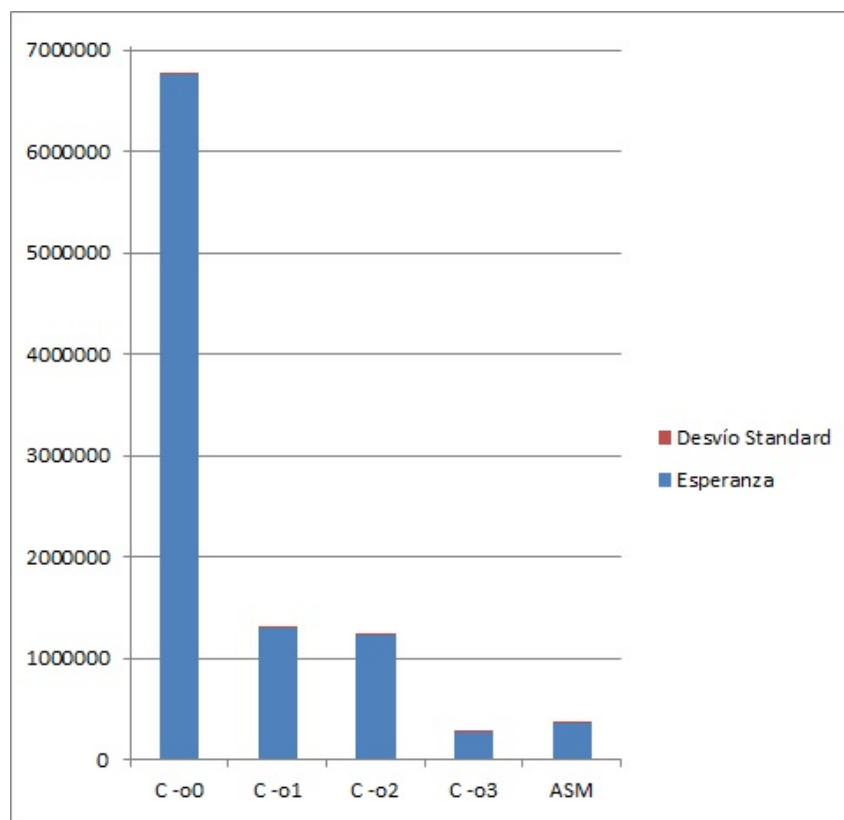
Experimento 1.4 - secuencial vs. vectorial

En este experimento deberá realizar una medición de las diferencias de performance entre las versiones de C y ASM (el primero con -O0, -O1, -O2 y -O3) y graficar los resultados.

El siguiente gráfico indica la esperanza de la cantidad de ciclos de clock que toma ejecutar el filtro Cropflip con los parámetros 404 404 4 4 en ASM y en C variando los flags de o0 a o3.

Reflejando las siguientes magnitudes:

	Esperanza	Desvío estándar
C -o0	6.761.044,5	131,463
C -o1	1.300.109	53,447
C -o2	1.227.102,5	50,143
C -o3	266.688	0,287
ASM	362.232	2,079



Se puede observar que la varianza es casi despreciable considerando el valor de la esperanza. Además, es notorio cómo el correr el programa con la orden de -o0 no efectúa ninguna optimización. También se puede observar que el código corrido en C bajo el comando de -o3 tiene una esperanza menor a la del código Assembler. Esto se debe a que el flag de optimización -o3 lo que hace es utilizar instrucciones SSE para procesar 4 pixels por iteración, lo cual le brinda un boost de velocidad mucho mayor. Además de haberse reducido la cantidad de instrucciones. Tal cual lo explicado anteriormente. Y por este motivo, es una buena optimización para el código.

Experimento 1.5 - cpu vs. bus de memoria

Se desea conocer cual es el mayor limitante a la performance de este filtro en su versión ASM.

¿Cuál es el factor que limita la performance en este caso? En caso de que el limitante fuera la intensidad de cómputo, entonces podrían agregarse instrucciones que realicen accesos a memoria extra y la performance casi no debería sufrir. La inversa puede aplicarse, si el limitante es la cantidad de accesos a memoria.⁴

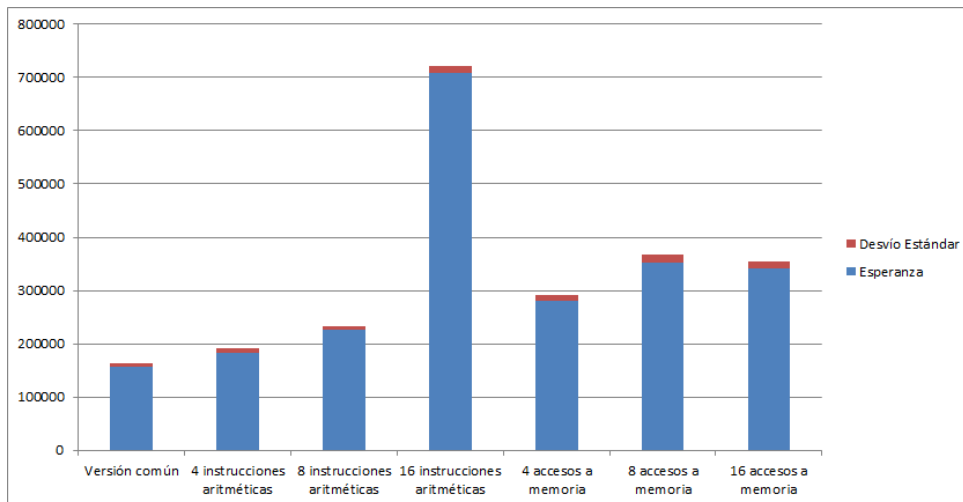
Realizar un experimento, agregando 4, 8 y 16 instrucciones aritméticas (por ej `add rax, rbx`) analizando como varía el tiempo de ejecución. Hacer lo mismo ahora con instrucciones de acceso a memoria, haciendo mitad lecturas y mitad escrituras (por ejemplo, agregando dos `mov rax, [rsp]` y dos `mov [rsp+8], rax`).⁵

Realizar un único gráfico que compare: 1. La versión original 2. Las versiones con más instrucciones aritméticas 3. Las versiones con más accesos a memoria

Acompañar al gráfico con una tabla que indique los valores graficados.

Las instrucciones aritméticas usadas son: `add eax,8` `sub eax 8` repetidas la cantidad de veces necesitada y los accesos a memoria realizados fueron siempre a la imagen fuente o la de destino, ya que tenemos la certeza de poseer el permiso para acceder a esas direcciones.

Cropflip	Esperanza	Desvío estándar
Versión común	157.236,897	5.334,414
Con 4 instrucciones aritméticas	183.968,425	8.391,506
Con 8 instrucciones aritméticas	225.991,244	7.729,428
Con 16 instrucciones aritméticas	707.857,067	12.756,322
Con 4 accesos a memoria	281.032,176	11.258,907
Con 8 accesos a memoria	352.370,777	14.353,789
Con 16 accesos a memoria	342.023,334	12.136,799



Como era de esperar, la versión original es la que menos tiempo demanda. Al agregarle 4 y 8 instrucciones aritméticas, se observa un incremento gradual que se condice con la cantidad de instrucciones. Al agregarle 16 instrucciones se puede observar un pico en la medición, el cual nos hace pensar que la

⁴también podría pasar que estén más bien balanceados y que agregar cualquier tipo de instrucción afecte sensiblemente la performance

⁵Notar que en el caso de acceder a `[rbp]` o `[rsp+8]` probablemente haya siempre hits en la cache, por lo que la medición no será de buena calidad. Si se le ocurre la manera, realizar accesos a otras direcciones alternativas.

forma en que se aumenta el tiempo de ejecución no es lineal respecto de la cantidad instrucciones, aunque no hay que omitir que de 4 a 8 y a 16 se va incrementando en potencia de dos (lo cual tampoco es lineal). En el caso de los accesos a memoria no sucede lo mismo que lo descrito anteriormente, cuando se agregan 4 accesos aumenta la cantidad de tiempo respecto de la original, también es así cuando se realizan 8 accesos, pero al realizar 16 accesos en nuestras mediciones disminuye el valor de su esperanza e incluso el de la varianza. Esto no es un indicador con el que podamos afirmar que siempre al realizar 16 accesos a memoria va a tardar menos que al hacer 8. Lo que intuimos que puede suceder es que, el aumento tampoco acá es lineal y para ese intervalo (8,16) el incremento es paulatino. Tampoco hay que dejar de lado que toda medición contiene cierto error y no es exacta y que, a pesar de haber eliminado los outliers, el resultado obtenido se aleje un poco de la “realidad”.

4.3. Filtro *Sierpinski*

Programar el filtro *Sierpinski* en lenguaje C y en en ASM haciendo uso de las instrucciones vectoriales (SSE).

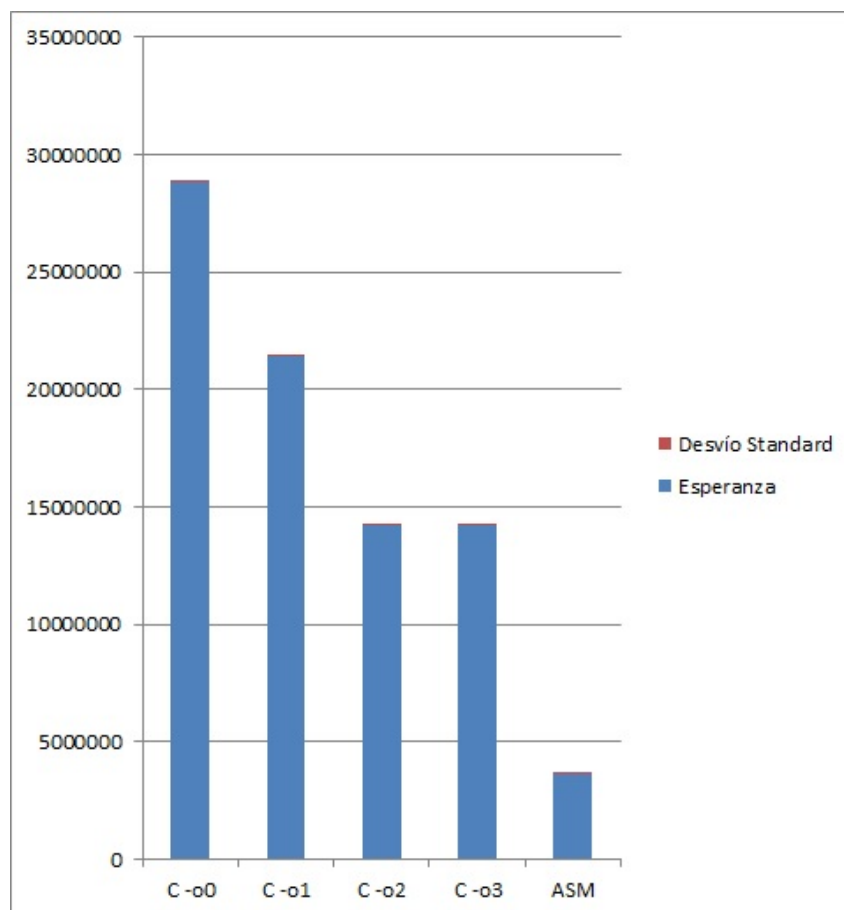
Experimento 2.1 - secuencial vs. vectorial

Analizar cuales son las diferencias de performace entre las versiones de C y ASM de este filtro, de igual modo que para el experimento 1.4.

El gráfico a continuación indica la esperanza y el desvío estándar de la cantidad de ciclos de clock que toma ejecutar el filtro *Sierpinski* en ASM y en C variando los flags de o0 a o3.

Reflejando las siguientes magnitudes:

	Esperanza	Desvío estándar
C -o0	28.801.586,5	78,016
C -o1	21.395.732	78,447
C -o2	14.231.758	116,445
C -o3	14.229.650	116,022
ASM	3.661.626	115,583



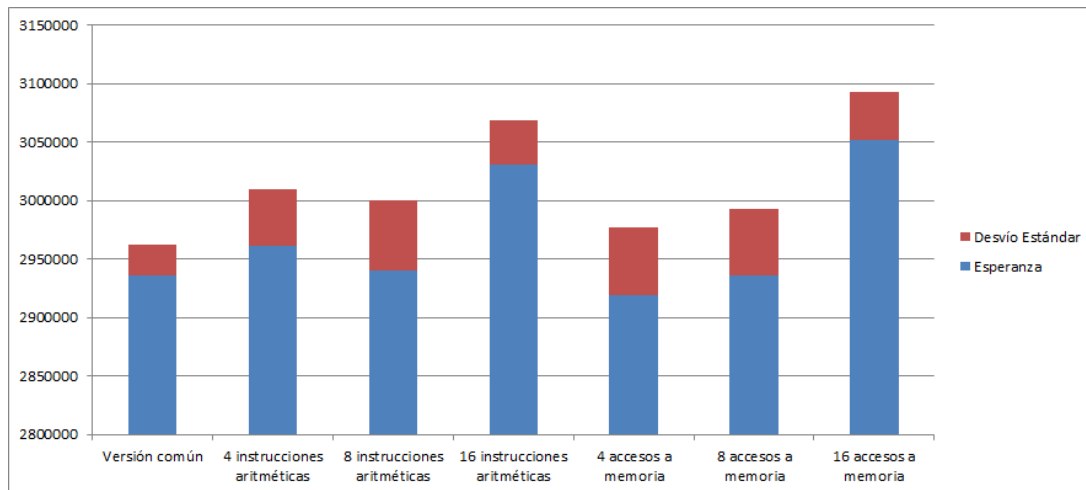
Aquí también el valor del desvío estándar es despreciable. Donde correr el código en C bajo el comando -o0 sigue siendo el caso con mayor esperanza (peor caso) y además la esperanza del código Assembler es menor a la del código C con -o3. Las optimizaciones -o2 y -o3 arrojan resultados similares, ambos por encima de los valores en ASM. En este caso, notablemente es más eficaz el filtro programado en Assembler.

Experimento 2.1 - cpu vs. bus de memoria

¿Cuál es el factor que limita la performance en este filtro? Repetir el experimento 1.5 para este filtro.

Las instrucciones utilizadas son las mismas que en el caso anterior, lo mismo con los accesos a memoria.

	Esperanza	Desvío estándar
Versión común	2.936.302,152	26.416,331
Con 4 instrucciones aritméticas	2.961.773,596	47.749,307
Con 8 instrucciones aritméticas	2.940.553,343	59.912,727
Con 16 instrucciones aritméticas	3.031.280,298	37.499,415
Con 4 accesos a memoria	2.919.253,679	57.697,386
Con 8 accesos a memoria	2.936.432,138	56.870,171
Con 16 accesos a memoria	3.051.926,505	40.964,812



En este caso, lo primero que nos parece necesario destacar es la magnitud del desvío estándar en todos los casos; son desvíos estándar grandes, dados los valores de la esperanza. Esto nos indica que los resultados de las mediciones obtenidas pueden no ser tan precisos como necesitamos ya que hay mucha dispersión.

Teniendo en cuenta lo mencionado arriba, se puede explicar porqué al sumarle al programa 4 instrucciones aritméticas aumenta la cantidad de tiempo necesaria respecto de la versión original, al igual que al sumarle 16 instrucciones pero no sucede lo mismo en el caso de las 8 instrucciones, ya que la esperanza es muy cercana numéricamente a la versión original. Observando el gráfico, queda expuesto notablemente que el desvío estándar de este caso es el mayor (siendo más del doble que en la versión común). Asumimos que en el caso “ideal” el crecimiento es paulatino y el tiempo necesario para ejecutar las 8 instrucciones queda por encima del caso anterior y por debajo del siguiente.

En el caso de los accesos a memoria, si bien también presentan un notoriamente elevado valor de desvío estándar para todos los casos, el comportamiento es más intuitivo. Se puede apreciar en el gráfico que el incremento del tiempo es acorde al aumento de los accesos a memoria y se da de manera pareja sin presentar ningún tipo de “pico”.

4.4. Filtro *Bandas*

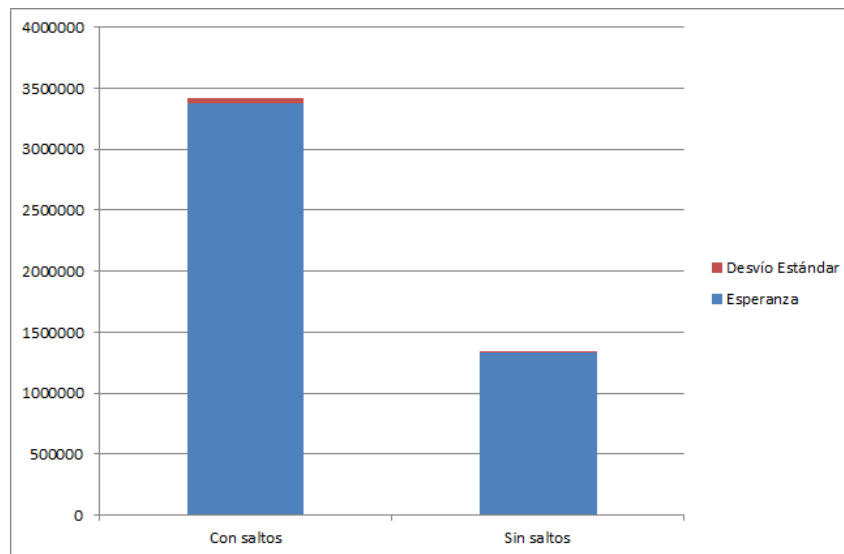
Programar el filtro *Bandas* en lenguaje C y en en ASM haciendo uso de las instrucciones vectoriales (SSE).

Experimento 3.1 - saltos condicionales

Se desea conocer que tanto impactan los saltos condicionales en el código de filtro *Bandas* con -01 (la versión en C).

Para poder medir esto de manera aproximada, remover el código que detecta a que banda pertenece cada pixel, dejando sólo una banda. Por más que la imagen resultante no sea correcta, será posible tomar una medida aproximada del impacto de los saltos condicionales. Analizar como varía la performance.

En la siguiente figura se ve cómo varía la esperanza y el desvío estándar entre dos corridas de C con el flag -o1 ambas, el Caso 1 es la versión original y el Caso 2 es comentando la guarda de las bandas:



Los valores del gráfico son:

Versión original: esperanza de 3.374.834,128 y varianza de 48.256,45868.

Versión sin guarda: esperanza de 1.333.624,786 y varianza de 2.545,912732.

En el gráfico anterior se puede ver que la influencia de los saltos condicionales es notable, ya que disminuye de manera notable la esperanza, siendo en ambos casos la varianza pequeña.

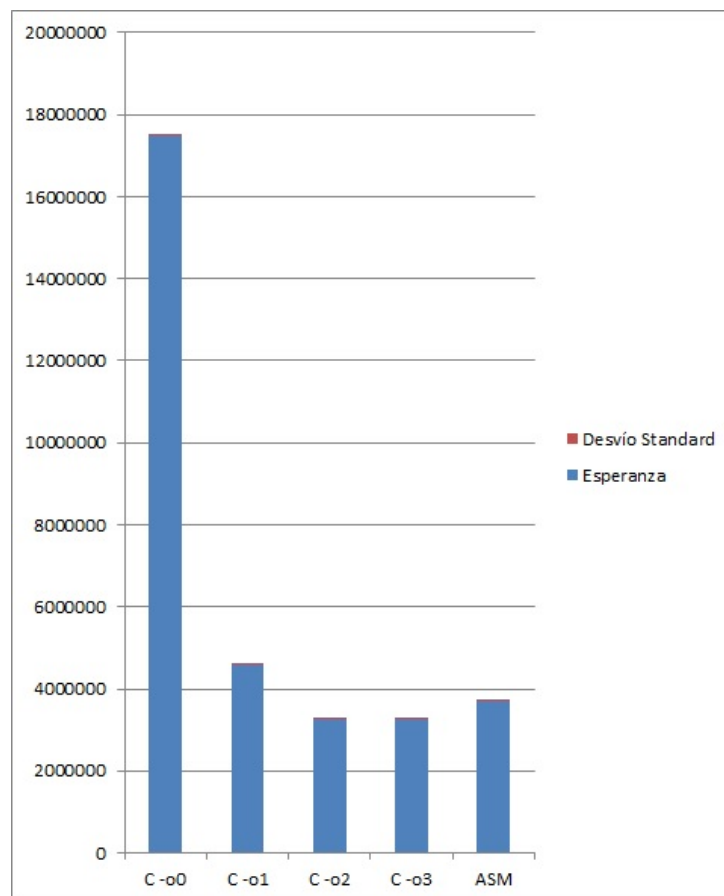
Experimento 3.2 - secuencial vs. vectorial

Repetir el experimento 1.4 para este filtro.

El siguiente gráfico indica la esperanza de la cantidad de ciclos de clock que toma ejecutar el filtro *Bandas* en ASM y en C variando los flags de o0 a o3.

Reflejando las siguientes magnitudes:

	Esperanza	Desvío estándar
C -o0	17.472.791	120,620
C -o1	4.583.340	128,514
C -o2	3.259.839	117,746
C -o3	3.259.767	117,391
ASM	3.703.203,5	115,659



Aquí también el valor del desvío estándar es despreciable. En este caso, podemos notar una vez más que el comando -o0 no ejerce ninguna optimización para este programa, ya que es el que tiene mayor tiempo de ejecución. Bajo el flag -o1 presenta una mejoría, pero sigue requiriendo un tiempo mayor que el programa en Lenguaje Assembler para su ejecución. Sin embargo, sucede algo llamativo que es que corriendo el lenguaje en lenguaje C tanto bajo el flag -o2 como para -o3 presentó una esperanza menor que el código en Assembler. Asumimos que ocurre esto porque estos flags están diseñados para casos como este, donde ensamblan un código fuente mejor que el nuestro en Assembler.

4.5. Filtro *Motion Blur*

Programar el filtro *mblur* en lenguaje C y en ASM haciendo uso de las instrucciones **SSE**.

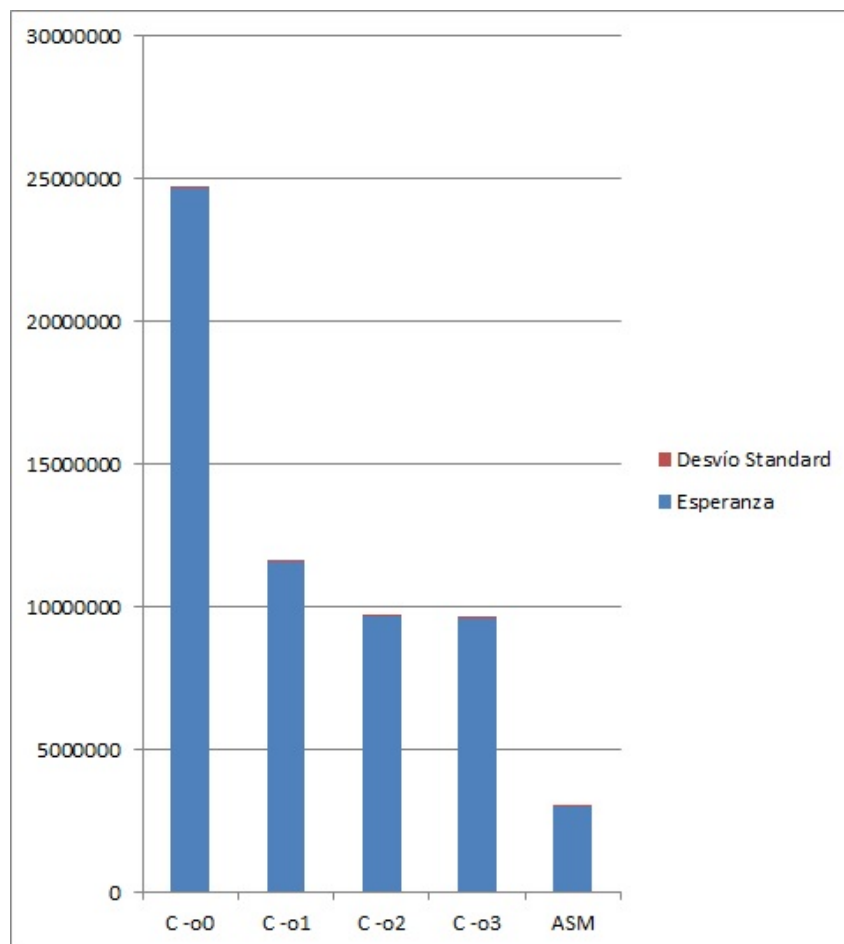
Experimento 4.1

Repetir el experimento 1.4 para este filtro

El siguiente gráfico indica la esperanza de la cantidad de ciclos de clock que toma ejecutar el filtro Bandas en ASM y en C variando los flags de o0 a o3.

Reflejando las siguientes magnitudes:

	Esperanza	Desvío estándar
C -o0	24.652.732,5	117,738
C -o1	11.562.261	103,024
C -o2	9.641.346	78,159
C -o3	9.639.625	78,227
ASM	2.993.980,5	118,245



Aquí también el valor del desvío estándar es despreciable. Donde las optimizaciones del compilador en la versión -o3 y la versión -o2 no son suficientes para tener una esperanza menor a la esperanza obtenida bajo el código Assembler.

5. Conclusiones y trabajo futuro

Los resultados de nuestras mediciones han demostrado que la implementación de los filtros en Lenguaje Assembler han logrado una performance de alto nivel de rendimiento ya que el tiempo necesario para su ejecución fue en la mayoría de los casos el menor. Si bien, hubo casos donde usando un flag de optimización adecuado al momento de compilar se lograron resultados con un tiempo de ejecución menor, no eran valores muy lejanos a los arrojados bajo el Lenguaje Assembler.

Dada la popularidad de la arquitectura de 64 bits de intel, se puede concluir que la implementación del SIMD en assembler tiene un precio razonable frente a la performance alcanzada.